This application is submitted in the name of the following inventor(s): 1 2 3 Residence City and State Citizenship **Inventor** Palo Alto, California 4 Blake LEWIS USA \$ Sunnyvale, California John EDWARDS **USA** 6 India Fremont, California Srinivasan VISWANATHAN 7 The assignee is Network Appliance, Inc., a California corporation 8 9 having an office at 495 East Java Drive, Sunnyvale, CA 94089. 10 Title of the Invention 11 12 13 **Instant Snapshot**

17 1. Field of Invention

14

15

16

18

This invention relates to data storage systems.

Background of the Invention

103.1035.01

This invention relates to data storage systems.

2

3

1

2. Related Art

4

5

6

7

8

9

10

11

12

13

Snapshots of a file system capture the contents of the files and directories in a file system at a particular point in time. Such snapshots have several uses. They allow the users of the file system to recover earlier versions of a file following an unintended deletion or modification. The contents of the snapshot can be copied to tape to provide a backup copy of the file system, and it can be copied to another file server and used as a replica. File systems, including the WAFL (Write Anywhere File Layout) file system, include a copy-on-write snapshot mechanism. Snapshot block ownership is recorded by updating the block's entry in the block map file, a bitmap associated with the vacancy of blocks.

OSELET TETEO 14

15

16

17

18

One problem with the prior art of creating snapshots is that the requirement for additional file system metadata in the active filesystem to keep track of which blocks snapshots occupy. This metadata requires 4 bytes per 4-KB file system block, i.e., 1/1024th of the file system. These methods are inefficient both in their use of storage space and in the time needed to create the snapshots.

A second problem with earlier snapshot implementations, was the time consuming steps of writing out a description of the snapshot state on creation and removing it on deletion.

A third problem with earlier copy-on-write mechanisms, was the required steps consumed a considerable amount of time and file system space. For example, some systems, such as those supplied with DCE/DFS include a copy-on-write mechanism for creating snapshots (called "clones"). The copy-on-write mechanism was used to record which blocks each clone occupied. Such systems require a new copy of the inode file and the indirect blocks for all files and directories are created when updating all of the original inodes.

Accordingly, it would be advantageous to provide an improved technique for more quickly and efficiently capturing the contents of the files and directories in the file system at a particular point in time. This is achieved in an embodiment of the invention that is not subject to the drawbacks of the related art.

Summary of the Invention

2

3

4

The invention provides an improved method and apparatus for creating a snapshot of a file system.

5

6

7

8

9

10

11

12

In a first aspect of the invention, a "copy-on-write" mechanism is used. An effective snapshot mechanism must be efficient both in its use of storage space and in the time needed to create it because file systems are often large. The snapshot uses the same blocks as the active file system until the active file system is modified. Whenever a modification occurs, the modified data is copied to a new block and the old data is saved (henceforth called "copy-on-write"). In this way, the snapshot only uses space where it differs from the active file system, and the amount of work required to create the snapshot initially is small.

13

14

15

16

17

In a second aspect of the invention, a record of which blocks are being used by the snapshot is included in the snapshot itself, allowing effectively instantaneous snapshot creation and deletion.

In a third aspect of the invention, the state of the active file system is described by a set of metafiles; in particular, a bitmap (henceforth the "active map") describes which blocks are free and which are in use by the active file system. The inode file describes which blocks are used by each file, including the metafiles. The inode file itself is described by a special root inode, also known as the "fsinfo block. This copy of the root inode becomes the root of the snapshot. The root inode captures all required states for creating the snapshot such as the location of all files and directories in the file system, it. During subsequent updates of the active file system, the system consults the bitmap included in the snapshot (the "snapmap") to determine whether a block is free for reuse or belongs to a snapshot. This mechanism allows the active file system to keep track of which blocks each snapshot uses without recording any additional bookkeeping information in the file system.

In a fourth aspect of the invention, a snapshot can also be deleted instantaneously simply by discarding its root inode. Further bookkeeping is not required, because the snapshot includes it's own description.

In a fifth aspect of the invention, the performance overhead associated with

the search for free blocks is reduced by the inclusion of a summary file. The summary file identifies blocks that are used by at least one snapshot; it is the logical OR of all the snapmap files. The write allocation code decides whether a block is free by examining the active map and the summary file. The active map indicates whether the block is currently in use in the active file system. The summary file indicates whether the block is used by any snapshot.

In a sixth aspect of the invention, the summary file is updated in the background after the creation or deletion of a snapshot. This occurs concurrently with other file system operations. Two bits are stored in the file system "fsinfo block" for each snapshot. These two bits indicate whether the summary file needs to be updated using the snapshot's snapmap information as a consequence of its creation or deletion. When a block is freed in the active file system, the corresponding block of the summary file is updated with the snapmap from the most recently created snapshot, if this has not already been done. An in-core bit map records the completed updates to avoid repeating them unnecessarily. This ensures that the combination of the active bitmap and the summary file will consistently identify all blocks that are currently in use. Additionally, the summary file is updated to reflect the effect of any recent snapshot deletions when

freeing a block in the active file system. This allows reuse of blocks that are now entirely

free. After updating the summary file following a snapshot creation or deletion, the

corresponding bit in the fsinfo block is adjusted.

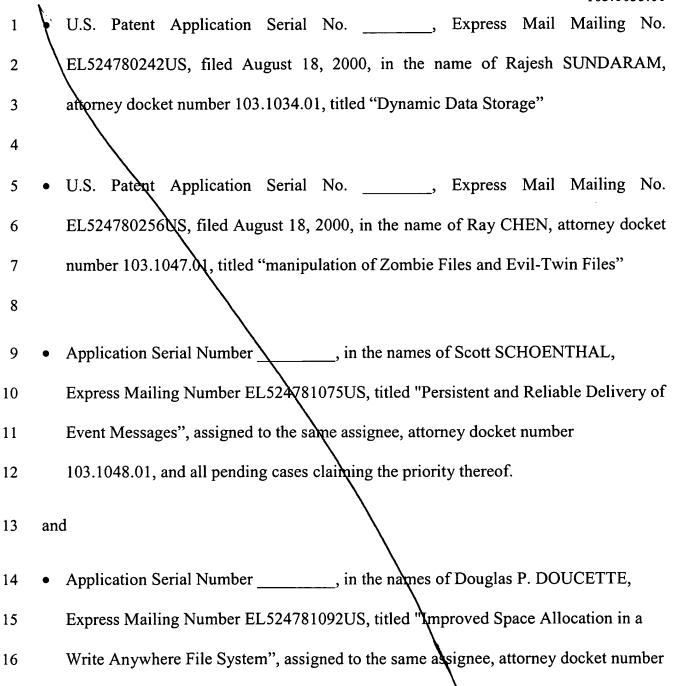
In a seventh aspect of the invention, the algorithm for deleting a snapshot involves examining the snapmaps of the deleted snapshot and the snapmaps of the next oldest and next youngest snapshot. A block that was used by the deleted snapshot but is not used by its neighbors can be marked free in the summary file, as no remaining snapshot is using it. However, these freed blocks cannot be reused immediately, as the snapmap of the deleted snapshot must be preserved until summary updating is complete. During a snapdelete free blocks are found by using the logical OR of the active bitmap, the summary file, and the snapmaps of all snapshots for which post-deletion updating is in progress. In other words, the snapmap of the deleted snapshot protects the snapshot

In the preferred embodiment, the invention is operative on WAFL file system. However, it is still possible for the invention to be applied to any computer data storage system such as a database system or a store and forward system such as cache or

from reuse until it is no longer needed for updating.

RAM if the data is kept for a limited period of time. 1 2 Brief Description of the Drawings 3 4 Figure 1 shows a block diagram of a system for an instant snapshot. 5 6 Figure 2 shows a block diagram of an instant snapshot. 7 8 Figure 3 shows a flow diagram of a method for creating a snapshot. 9 10 **Incorporated Disclosures** 11 12 The inventions described herein can be used in conjunction with inventions 13 described in the following applications: 14 **NS** · 15 U.S. Patent Application Serial No. , Express Mail Mailing No. EL 524781089US, filed August 18, 2000, in the name of Blake LEWIS, attorney docket 17 number 103,1033.01, titled "Reserving File System Blocks" 18 19

ÚSCLECE.



103.1045.01, and all pending cases claiming the priority thereof.

2

Detailed Description of the Preferred Embodiment

3 4

5

6

7

8

9

10

12

13

14

In the following description, a preferred embodiment of the invention is described with regard to preferred process steps and data structures. However, those skilled in the art would recognize, after perusal of this application, that embodiments of the invention might be implemented using a variety of other techniques without undue experimentation or further invention, and that such other techniques would be within the scope and spirit of the invention.

Lexicography

11

"CETEQU

As used herein, use of the following terms refer or relate to aspects of the invention as described below. The general meaning of these terms is intended to be illustory and in no way limiting.

15

16

17

18

Inode - In general, the term "inode" refers to data structures that include information about files in Unix and other file systems. Each file has an inode and is identified by an inode number (i-number) in the file system where it resides. Inodes provide

- important information on files such as user and group ownership, access mode (read,
- write, execute permissions) and type. An inode points to the file blocks or indirect
- 3 blocks of the file it represents.

- Sector In general, the term "sector" refers to a physical section of a disk drive
- 6 including a collection of bytes, such as 512 bytes.

7

- Data Storage Block In general, the phrase "data storage block" refers to specific
- 9 allocation areas on a hard disk. The allocation area is a collection of sectors, such as
- 8 sectors or 4,096 bytes, commonly called 4K bytes or 4-KB.

11

12

- File Block In general, the phrase "file block" refers to a standard size block of data
- including some or all of the data in a file. In the preferred embodiment, the file block
- is the same size as a data storage block.

- fsinfo (File System Information Block) In general, the phrase "file system
- information block" refers to one or more copies of a block known as the "fsinfo
- block". These blocks are located at fixed locations on the disks. The fsinfo block

includes data about the volume including the size of the volume, volume level

2 options, language and more.

3

• WAFL (Write Anywhere File Layout) - In general, the term "WAFL" refers to a high level structure for a file system. Pointers are used for locating data. All the data is included in files. These files can be written anywhere on the disk in chunks of file

7 blocks placed in data storage blocks.

8

9

• Volume In general, the term "volume" refers to a single file system. The file system may be composed of a collection of disk drives.

11

12

14

15

16

17

10

Consistency Point (CP) - In general, the term "CP" refers to a time that a file system reaches a consistent state. When this state is reached, all the files have been written to all the blocks and are safely on disk and the one or more copies of redundant fsinfo blocks get written out. If the system crashes before the fsinfo blocks go out, all other changes are lost and the system reverts back to the last CP. The file system advances atomically from one CP to the next.

• Consistent State - In general, the phrase "consistent state" refers to the system configuration of files in blocks after the CP is reached.

3

• Range - In general, the term "range" refers to a group of blocks, such as 1,024 blocks.

5

• Active file system - In general, the phrase "active file system" refers to the current

file system arrived at with the most recent CP. In the preferred embodiment, the

active file system includes the active map, the summary map and points to all

snapshots and other data storage blocks through a hierarchy of inodes, indirect data

storage blocks and more.

11

12

13

• Active map - In general, the phrase "active map" refers to a to a file including a bitmap associated with the vacancy of blocks of the active file system.

14

• Snapshot - In general, the term "snapshot" refers to a copy of the file system. The snapshot diverges from the active file system over time as the active file system is modified. A snapshot can be used to return the file system to a particular CP (consistency point).

8

9

10

11

- Snapmap In general, the term "snapmap" refers to a file including a bitmap
 associated with the vacancy of blocks of a snapshot. The active map diverges from a
 snapmap over time as the blocks used by the active file system change during
 consistency points.
- Summary map In general, the term "summary map" refers to a file including an IOR (inclusive OR) bitmap of all the snapmaps.
 - Space map In general, the term "space map" refers to a file including an array of numbers which describe the number of storage blocks used in an allocation area.
- Blockmap In general, the term "blockmap" refers to a map describing the status of the blocks in the file system.
- Snapdelete In general, the term "snapdelete" refers to an operation that removes a

 particular snapshot from the file system. This command can allow a storage block to

 be freed for reallocation provided no other snapshot or the active file system uses the

 storage block.

Snapcreate – In general, the term "snapcreate" refers to the operation of retaining a
 consistency point and preserving it as a snapshot.

3

4

5

6

As described herein, the scope and spirit of the invention is not limited to any of the definitions or specific examples shown therein, but is intended to include the most general concepts embodied by these and other terms.

7

8 System Elements

9

Figure 1 shows a block diagram of a system for an instant snapshot.

11

12

13

14

15

10

The root block 100 includes the inode of the inode file 105 plus other information regarding the active file system 110, the active map 115, previous active file systems known as snapshots 120, 125, 130 and 135 and their respective snapmaps 140, 145, 150 and 155.

16

17 The active map 115 of the active file system 110 is a bitmap associated with the vacancy of blocks for the active file system 110. The respective snapmaps 140,

1 145, 50 and 155 are active maps can be associated with particular snapshots 120, 125,

2 130 and 35 and an inclusive OR summary map 160 of the snapmaps 140, 145, 150 and

3 155. Also shown are other blocks 115 including double indirect blocks 130 and 132,

4 indirect blocks 165, 166 and 167 and data blocks 170, 171, 172 and 173. Finally, Figure

1 shows the spacemap 180 including a collection of spacemap blocks of numbers 181,

6 182, 183, 184 and 190.

7

8

9

10

11

12

5

The root block 100 includes a collection of pointers that are written to the file system when the system has reached a new CP (consistency point). The pointers are aimed at a set of indirect (or triple indirect, or double indirect) inode blocks (not shown) or directly to the inode file 105 consisting of a set of blocks known as inode blocks 191, 192, 193, 194 and 195.

13

14

15

16

17

18

The number of total blocks determines the number of indirect layers of blocks in the file system. The root block 100 includes a standard quantity of data, such as 128 bytes. 64 of these 128 bytes describe file size and other properties; the remaining 64 bytes are a collection of pointers to the inode blocks 191, 192, 193, 194 and 195 in the inode file 105. Each pointer in the preferred embodiment is made of 4 bytes. Thus, there

are approximately 16 pointer entries in the root block 100 aimed at 16 corresponding

inode blocks of the inode file 105 each including 4K bytes. If there are more than 16

inode blocks, indirect inode blocks are used.

In a preferred embodiment, file blocks are 4096 bytes and inodes are 128 bytes. It follows that each block of the inode file contains 32 (i.e. 4,096/128) separate inodes that point to other blocks 115 in the active file system.

Inode block 193 in the inode file 105 points to a set of blocks (1, 2, 3, ..., P) called the active map 115. Each block in the active map 115 is a bitmap where each bit corresponds to a block in the entire volume. A "1" in a particular position in the bitmap correlates with a particular allocated block in the active file system 110. Conversely, a "0" correlates to the particular block being unused by the active file system 110. Since each block in the active map 115 can describe up to 32K blocks or 128 MB, 8 blocks are required per GB, 8K blocks per TB.

Another inode block in the inode file 105 is inode block N 212. This block includes a set of pointers to a collection of snapshots 120, 125, 130 and 135 of the

volume. Each snapshot includes all the information of a root block and is equivalent to an older root block from a previous active file system. The snapshot 120 may be created 2 at any past CP. Regardless when the snapshot is created, the snapshot is an exact copy 3 of the active file system at that time. The newest snapshot 120 includes a collection of 4 pointers that are aimed directly or indirectly to the same inode file 105 as the root block 5 100 of the active file system 110. As the active file system 110 changes (generally from 6 writing files, deleting files, changing attributes of files, renaming file, modifying their 7 contents and related activities), the active file system and snapshot will diverge over time. 8 Given the slow rate of divergence of an active file system from a snapshot, any two 9 snapshots will share many of the same blocks. The newest snapshot 120 is associated 10 with snapmap 140. Snapmap 140 is a bit map that is initially identical to the active map 11 115. The older snapshots 125, 130 and 194 have a corresponding collection of snapmaps 12 145, 150 and 155. Like the active map \(\)\(\)5, these snapmaps 145, 150 and 155 include a 13 set of blocks including bitmaps that correspond to allocated and free blocks for the 14 particular CP when the particular snapmaps 145, 150 and 155 were created. Any active 15 file system may have a structure that includes pointers to one or more snapshots. 16 Snapshots are identical to the active file system when they are created. It follows that 17 snapshots contain pointers to older snapshots. There can be a large number of previous 18

snapshots in any active file system or snapshot. In the event that there are no snapshot,

2 there will be no pointers in the active file system.

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

1

Blocks not used in the active file system 110 are not necessarily available for allocation or reallocation because the blocks may be used by snapshots. Blocks used by snapshots are freed by removing a snapshot using the snapdelete command. When a snapshot is deleted any block used only by that snapshot and not by other snapshots nor by the active file system becomes free for reuse by WAFL. If no other snapshot or active files uses the block, then the block can be freed and written over during the next copy onwrote-execution by WAFL. The system can relatively efficiently determine whether a block can be removed using the "nearest neighbor rule". If the previous and next snapshot do not allocate a particular block in their respective snapmaps, then the block can be freed for reuse by WAFL. For WAFL to find free space to write new data or metadata, it could search the active map 115 and the snapmaps (140, 145, 150 and 155) of the snapshots (120, 125, 130 and 135) to find blocks that are totally unused. This would be very inefficient; thus it is preferable to use the active map and the summary map as described below

15

16

17

18

7

9

A5 7

A summary map 160 is created by using an IOR (inclusive OR) operation 139 on the snapmaps 140, 145, 150 and 155. Like the active map 115 and the snapmaps 2 140, 145, 150 and 155, the summary map 160 is a file whose data blocks (1, 2, 3, ...Q) 3 contained a bit map. Each bit in each block of the summary map. describes the allocation 4 status of one block in the system with "1" being allocated and "0" being free. The 5 summary map 160 describes the allocated and free blocks of the entire volume from all 6

the snapshots 120, 125, 130 and 135 combined. The use of the summary file 160 is to

8 avoid overwriting blocks in use by snapshots?

An IOR operation on sets of blocks (such as 1,024 blocks) of the active map 115 and the summary map 160 produces a spacemap 180. Unlike the active map 115 and the summary map 160, which are a set of blocks containing bitmaps, the spacemap 180 is a set of blocks including 181, 182, 183, 184, and 190 containing arrays of binary numbers. The binary numbers in the array represent the addition of all the vacant blocks in a region containing a fixed number of blocks, such as 1,024 blocks. The array of binary numbers in the single spacemap block 181 represents the allocation of all blocks for all snapshots and the active file system in one range of 1,024 blocks. Each of the binary numbers 181, 182, 183, 184, and 190 in the array are a fixed length. In a

preferred embodiment, the binary numbers are 16 bit numbers, although only 10 bits are

2 used.

3

4

5

6

7

8

9

10

11

12

1

In a preferred embodiment, the large spacemap array binary number 182 (0000001111111110=1,02) in decimal units) tells the file system that the corresponding In such embodiments, the largest binary number range is relatively full. 00001111111111 (1,023 in decimal) represents a range containing at most one empty... The small binary number 184 (000000000001110=13 in decimal units) instructs the file system that the related range is relatively empty. The spacemap 180 is thus a representation in a very compact form of the allocation of all the blocks in the volume broken into 1,024 block sections. Each 16 bit number in the array of the spacemap 180 corresponds to the allocations of blocks in the range containing 1,024 blocks or about 4 MB. Each spacemap block 180 has about 2,000 binary numbers in the array and they describe the allocation status for 8 GB. Unlike the summary map 120, the spacemap

16

15

Figure 2 shows a block diagram of an instant snapshot.

block 180 needs to be determined whenever a file needs to be written.

18

The old root block 200 of snapshot #1 201 includes the inode of the inode file 202 plus other information regarding the previous active file system known as snapshot #1 201, the snap map 205, earlier active file systems known as snapshot #2 210,

snapshot #3 215 and snapshot #4 220, and their respective snapmaps 225, 230 and 235.

5

7

8

10

11

12

13

4

bitmap associated with the vacancy of blocks for snapshot #1 201. The respective snapmaps 225, 230 and 235 are earlier active maps that can be associated with particular snapshots 210, 215 and 220 and an inclusive OR summary map 245 of the snapmaps 225, 230 and 235. Also shown are other blocks 211 including double indirect blocks 240 and 332, indirect blocks 250, 251 and 252 and data blocks 260, 262, 263 and 264. Finally, Figure 2 shows the spacemap 270 of snapshot #1 201 including a collection of spacemap blocks of binary numbers 272, 273, 274, 275 and 276.

14

15

16

17

18

The old root block 200 includes a collection of pointers that are written to the previous active file system when the system had reached the previous CP. The pointers are aimed at a set of indirect (or triple indirect, or double indirect) inode blocks (not shown) or directly to the inode file 202 consisting of a set of blocks known as inode

1 blocks 281, 282, 283, 284 and 285.

old root block 201 starting with double indirect blocks 240 and 332 (there could also be triple indirect blocks). The double indirect blocks 240 and 332 include pointers to indirect blocks 250, 251 and 252. The indirect blocks 250, 251 and 252 include pointers that are directed to data leaf blocks 260, 262, 263 and 264 of the active file system 201.

Inode block 283 in the inode file 202 points to a set of blocks (1, 2, 3, ..., P) called the snap map 205. Each block in the snap map 205 is a bitmap where each bit corresponds to a block in the entire volume. A "1" in a particular position in the bitmap correlates with a particular allocated block in the active file system 201. Conversely, a "0" correlates to the particular block being free for allocation in the old root block 201. Each block in the snap map 205 can describe up to 32K blocks or 128 MB.

Inode file 202 also includes inode block N 285. This block includes a set of pointers to a collection of earlier snapshots, snapshot #2 210, shapshot #3 215 and snapshot #4 220 of the volume. Each snapshot includes all the information of a root

1 block and is equivalent to an older root block from a previous active file system.

2

Spapshot #1 201 also includes an old summary map 245 and old spacemap

- blocks 270. Although these blocks of data are included in snapshot #1 201 and previous
- snapshots, in a preferred embodiment, this data is not used by the active file system of 5
- figure 2. 6

7

Method of Use 8

9

10

Figure 3 shows a flow diagram of a method for using a system as shown in

figure 1. 11

COSTROS. Lacrado 12

method 300 is performed by the file system 100. Although the method

- 400 is described serially, the steps of the method 300 can be performed by separate
- elements in conjunction or in parallel, whether asynchronously, in a pipelined manner, or 15
- otherwise. There is no particular requirement that the method 300 be performed in the 16
- same order in which this description lists the steps, except where so indicated. 17

	4
	5
	6
	7
	8
	9
Į.	10
i i	11
	12
	13
	14
metra	15
	16

18

1	At a flow point 305, the file system 100 is ready to perform a method 300.
2	
3	At a step 310, a user will request a snapshot of the file system 100.
4	
5	At a step 315, a timer associated with the file system 100 initiates the
6	creation of a new snapshot.
7	
8	At a step 320, the file system 100 receives a request to make a snapshot.
9	
10	At a step 325, the file system 100 creates a new file.
11	
12	At a step 330, the root node of the new file points to the root node of the
13	current active file system.
14	
15	At a step 335, the file system 100 makes the file read only.

At a step 340, the file system 100 updates the new summary map by using

an inclusive OR of the most recent snapmap and the existing summary file. This step

	8
<u>0</u>	9
FNO.	10
) -	11
	12
H	13
	14

1	must be	done	before	any	blocks	are	freed	in	the	corresponding	active	map	block.	If
---	---------	------	--------	-----	--------	-----	-------	----	-----	---------------	--------	-----	--------	----

- multiple snapshots are created such that the processing overlaps in time, the update in 2
- step 340 need only be done for the most recently created snapshot. 3

5

6

At a flow point 345, the snapshot create and the summary file update is completed and the snapshot creation is done.

7

An analogous method may be performed for snapshot delete.

Alternative Embodiments

15

Although preferred embodiments are disclosed herein, many variations are possible which remain within the concept, scope, and spirit of the invention, and these variations would become clear to those skilled in the art after perusal of this application.